



# Abstracts Embeddings Evaluation

A Case Study of Artificial Intelligence and Medical Imaging for the  
COVID-19 Infection

**Giovanni Zurlo, Elisabetta Ronchieri**

September 18, 2023

ML\_INFNO, Bologna, IT





# Table of Contents

## 1 Introduction

- ▶ Introduction
- ▶ Methodology
- ▶ Results
- ▶ Conclusions and Future Work



# Problem Domain

## 1 Introduction

- The SARS-CoV-2 pandemic triggered unprecedented research efforts across various disciplines.
- This study delves into the collaborative prospects of artificial intelligence (AI) and medical imaging to expedite the analysis of scientific COVID-19 articles on larger scale.
- By harnessing the capabilities of natural language processing (NLP) and contextualized vector representations, the investigation scrutinizes the potential of popular biomedical transformer-based models to capture the semantic attributes in the medical imaging literature.



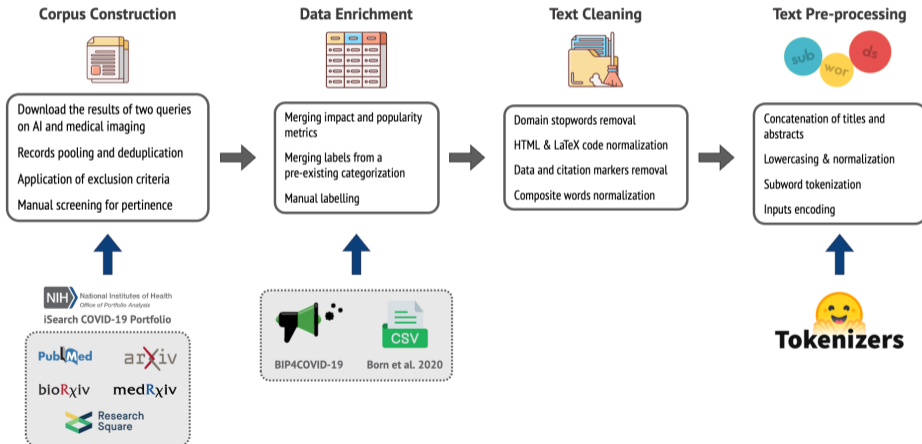
# Table of Contents

## 2 Methodology

- ▶ Introduction
- ▶ **Methodology**
- ▶ Results
- ▶ Conclusions and Future Work

# Data Collection Workflow

## 2 Methodology





# Papers Sources

## 2 Methodology

### Considered the National Institutes of Health (NIH)'s iSearch COVID-19 Portfolio

NIH National Institutes of Health Office of Portfolio Analysis COVID-19

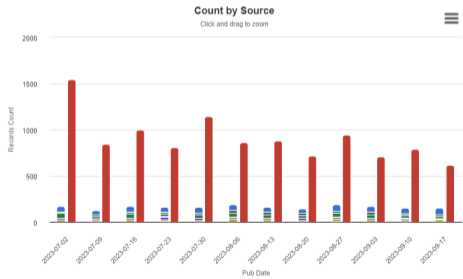
Search Query

AND OR

Q - X T

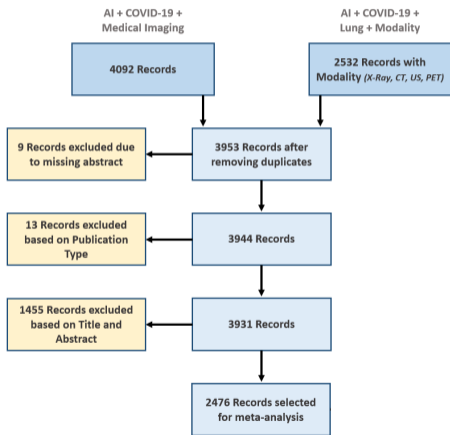
Welcome to the COVID-19 Portfolio

The *iSearch* COVID-19 portfolio is NIH's comprehensive, expert-curated source for publications and preprints related to either COVID-19 or the novel coronavirus SARS-CoV-2. Our COVID-19 Portfolio tool leverages the cutting-edge analytical capability of the *iSearch* platform, with its powerful search functionality and faceting, and includes



# PRISMA-based Corpus Definition

## 2 Methodology



### 1. Broad query

AI AND COVID-19 AND 'Medical Imaging'

### 2. Modality-specific query

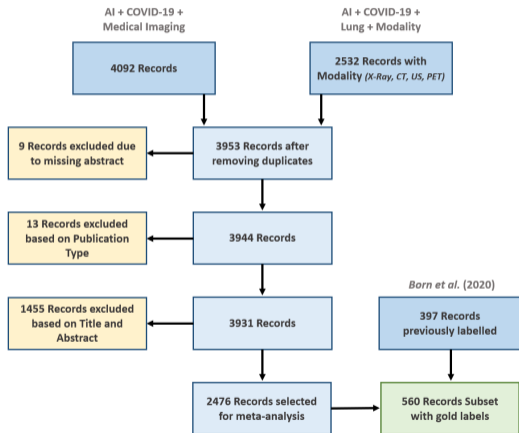
AI AND COVID-19 AND Lung AND (CT OR CXR OR US OR PET)

- CT Computerized Tomography
- CXR Chest X-Ray
- US Ultrasound
- PET Positron Emission Tomography

Collected papers  $\in$  period Jan 1, 2020 - May 27, 2023  
Used Preferred Reporting Items for Systematic Reviews and Meta-Analysis (PRISMA) to select papers

# Data Enrichment

## 2 Methodology



## Identified 560 gold-papers

- 163 papers manually labeled to address the issue of class imbalance
- 397 papers derived by Born et al. (2020) (see ref [6] in the paper) based on a supplementary dataset titled 'Detailed results of systematic meta-analysis', merged by using title and already labeled





# Data Enrichment - Labeling Assignment

## 2 Methodology

Chose to adopt the tasks and modalities classification framework already adopted in Born et al. (2020) (see ref [6] in the paper)

### Primary Task

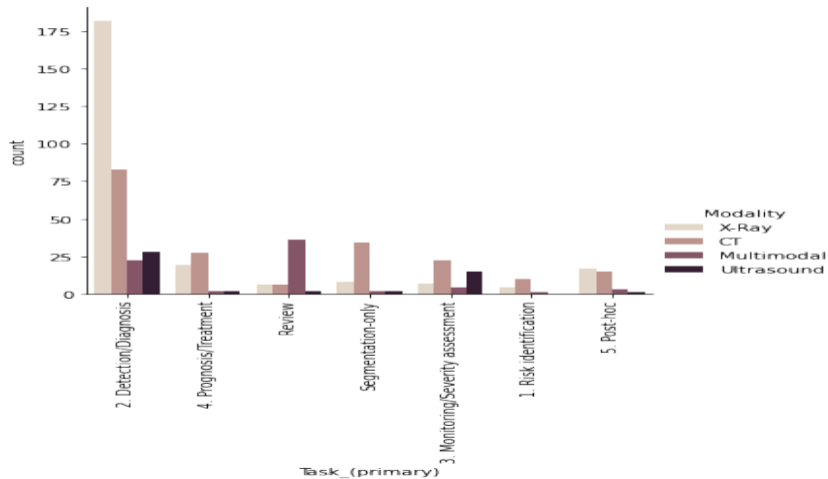
- Detection/Diagnosis
- Monitoring/Severity
- Assessment
- Post-Hoc
- Prognosis/Treatment
- Review
- Risk Identification
- Segmentation-only (for lung tissue or other disease features without any clinically relevant downstream tasks)

### Imaging Modality

- CT
- CXR
- Lung US
- Multimodal

# Data Enrichment - Labeling Assignment

## 2 Methodology





# 15 Models

## 2 Methodology

### 12 Bidirectional Encoder Representation (BERT) models

- original BERT in its base and large versions
- SciBERT
- BioBERT in its base and large versions
- PubMedBERT in its base and large versions
- COVID-19 BERT
- COVID SciBERT
- ClinicalCovidBERT
- RadBERT
- BioCovidBERT

### 3 SPECTER models

- standard model
- two-others with task-specific adapters



# Models Distinctions Summary

## 2 Methodology

Model	Training Corpus	Weights Initialization	Vocabulary	Details on Training Corpus
<b>BERT<sub>base</sub></b>	Wiki+Books	From-scratch	Derived from corpus	800M + 2.5B words, 1M steps
<b>SciBERT</b>	SemanticScholar Full-Texts	BERT <sub>base</sub>	Derived from corpus	1.14M Full-Texts, 18% from computer science and 82% from broad biomedical domain
<b>BioBERT<sub>base</sub></b>	PubMed abstracts	BERT <sub>Base</sub>	Same as BERT <sub>Base</sub>	Updated 2019. 4.5B Words, 1M steps
<b>PubMedBERT<sub>base</sub></b>	PubMed abstracts + PMC Full-Texts	From-scratch	Derived from corpus	Updated Feb. 2020. 16.8B Words, 100K steps
<b>CORD-19 BERT</b>	CORD-19 dataset	BERT <sub>Base</sub>	Same as BERT <sub>Base</sub>	Updated Early 2020
<b>CovidSciBERT</b>	CORD-19 dataset	SciBERT	Extended from SciBERT	Updated Early 2020
<b>ClinicalCovidBERT</b>	CORD-19 dataset	Bio+Clinical BERT [1]	Same as BERT <sub>Base</sub>	Full-Texts updated June 2020, 150K steps
<b>RadBERT</b>	Radiology Reports	BioBERT <sub>base</sub>	Same as BERT <sub>Base</sub>	4M reports from 600K unique patients treated at Stanford Health Care from 1992 to 2014
<b>SPECTER 2</b>	6M Triplets of Papers Citations	SciBERT	Same as SciBERT	Extended version of the <code>cite_prediction</code> dataset from [31]
<b>BERT<sub>large</sub></b>	Wiki+Books	From-scratch	Derived from corpus	800M + 2.5B words, 1M steps
<b>BioBERT<sub>large</sub></b>	PubMed abstracts	BERT <sub>Large</sub>	Derived from corpus	Updated 2019. 4.5B Words, 1M steps
<b>PubMedBERT<sub>large</sub></b>	PubMed abstracts	From-scratch	Derived from corpus	Updated Feb. 2020. 3.2B Words, 100K steps
<b>BioCovidBERT</b>	CORD-19 dataset	BioBERT <sub>large</sub>	Same as BioBERT <sub>large</sub>	Full-Texts updated June 2020, 200K steps



# BERTbase example

## 2 Methodology

- Each title + abstract pair gets concatenated.
- Any BERT<sub>base</sub> model encodes each pair in 768-dimensional latent space.
- The first token is a special classification token [CLS].
- The separator token [SEP] marks its end and separates titles from abstracts.
- The context window is fixed at 512 tokens (almost 300-400 words), causing a truncation for longer inputs.
  - Our dataset adheres this constraint of the context window.
  - Few records required truncation.



# Extraction Strategies

## 2 Methodology

- Our goal is to obtain a singular vector representation for each Text + Abstract, no one for each token.
- Three extraction strategies considered:
  - the first two involve extracting the final hidden state representation of the [CLS] token and the trailing [SEP] token
  - the third uses the mean-pooling strategy based on the second-to-last hidden states.



# Performance Metrics

## 2 Methodology

- Accuracy is computed for a k-Nearest Neighbors (kNN) classifier to provide an evaluation of embedding quality.
- All kNN-based metrics involved k=6 or k=13 exact nearest neighbors.
  - 'KneihborsClassifier' class from Scikit-Learn 1.2.2
  - All parameters were chosen performing a cross-validated grid search:
    - algorithm='auto', weights='distance', distance='cosine'
- To predict each test paper's label, kNN takes a weighted majority vote among the paper's NNs' labels in the training set.
- Neibhbors are weighted by the inverse of their cosine distance.



# Accuracy Details

## 2 Methodology

- For the accuracy, cross-validated values were averaged over the same 10-fold split.
- Additionally, a balanced version of accuracy was computed.
- The chance-level accuracy was calculated by using 'DummyClassifier' with strategy='stratified' to ignore the input features.





# Table of Contents

3 Results

- ▶ Introduction
- ▶ Methodology
- ▶ **Results**
- ▶ Conclusions and Future Work

# Quality metrics for the embeddings (Imaging Modality Prediction)

## 3 Results

10-fold kNN classification accuracy and balanced accuracy.

Hyperparameters:  $k = 13$ ,  $weights = distance$ ,  $distance = cosine$ .

Model	Accuracy (%)			Balanced Accuracy (%)		
	[CLS]	[SEP]	AVG	[CLS]	[SEP]	AVG
<b>BERT<sub>base</sub></b>	54.3	57.7	61.3	38.6	43.3	49.9
<b>SciBERT</b>	58.2	56.4	63.6	43.9	41	48.2
<b>BioBERT<sub>base</sub></b>	53.2	65.2	61.1	36	52.6	46.2
<b>PubMedBERT<sub>base</sub></b>	57.7	74.5	64.3	42.9	58.6	50.1
<b>CORD-19 BERT</b>	56.4	52.9	60.2	43.0	37.6	44.9
<b>CovidSciBERT</b>	64.5	60.5	62.7	49.9	47.9	50.6
<b>ClinicalCovidBERT</b>	65.4	64.5	63.4	53.2	50.8	49.6
<b>RadBERT</b>	57.9	57.9	58.2	38	38	39.7
<b>SPECTER 2</b>	<b>82.5</b>	<b>83.8</b>	68.8	<b>75.3</b>	<b>76.7</b>	57.8
<b>BERT<sub>large</sub></b>	50	58.8	60.2	34.7	43.4	44.8
<b>BioBERT<sub>large</sub></b>	54.4	62.1	65.5	39.8	47.2	51.3
<b>PubMedBERT<sub>large</sub></b>	57	61.3	60.4	40.3	44.2	45.3
<b>BioCovidBERT</b>	69.5	64.8	<b>69.6</b>	53.8	49	<b>58.3</b>
<b>Chance Level</b>	35.5 ±13			24.8 ±9		

SPECTER employs the [CLS] token, but we also applied the others for consistency.

BioCovidBERT<sub>large</sub> slightly outperformed with AVG pooling strategy due to its continual pre-trained on a COVID-19-based corpus.

# Quality metrics for the embeddings (Task Prediction)

## 3 Results

10-fold kNN classification accuracy and balanced accuracy.

Hyperparameters:  $k = 6$ ,  $weights = distance$ ,  $distance = cosine$ .

Model	Accuracy (%)			Balanced Accuracy (%)		
	[CLS]	[SEP]	AVG	[CLS]	[SEP]	AVG
<b>BERT<sub>base</sub></b>	59.6	58.2	64.8	27	28.4	33.9
<b>SciBERT</b>	62.7	63	68.6	33.3	31.5	38.3
<b>BioBERT<sub>base</sub></b>	63.9	70.2	69.6	28.6	40.5	40.2
<b>PubMedBERT<sub>base</sub></b>	66.8	70.7	67.5	34.7	42.3	36.9
<b>CORD-19 BERT</b>	65.0	60.7	65.9	33.7	25.9	34.4
<b>CovidSciBERT</b>	70.2	70.2	71.8	42.3	42.4	45.1
<b>ClinicalCovid BERT</b>	70.9	71.3	70	43.6	46.7	40.9
<b>RadBERT</b>	60.9	60.9	61.2	26.5	26.5	26.4
<b>SPECTER 2</b>	<b>75.4</b>	<b>74.5</b>	<b>74.1</b>	<b>56.6</b>	<b>55.9</b>	<b>51.5</b>
<b>BERT<sub>large</sub></b>	60	64.1	66.6	26.2	34	38.5
<b>BioBERT<sub>large</sub></b>	62	68.6	67.3	28.7	37.7	36
<b>PubMedBERT<sub>large</sub></b>	63	67.7	68.9	30	36.1	38.4
<b>BioCovidBERT</b>	66.6	68.2	70.9	37	39.5	42.5
<b>Chance Level</b>	36.1 ±6			14.8 ±9		

The Balanced accuracy scores decreased due to the presence of stronger class imbalance and lower recall values for 'post-hoc' and 'risk identification' classes.



# Table of Contents

## 4 Conclusions and Future Work

- ▶ Introduction
- ▶ Methodology
- ▶ Results
- ▶ Conclusions and Future Work



# Conclusions

## 4 Conclusions and Future Work

- A first version of a medical imaging dataset for the COVID-19 infection has been defined by following the PRISMA procedure in order to evaluate embeddings quality for abstracts texts.
  - Extrinsic evaluation fails if the embeddings are trained to serve in a wide range of different tasks.
- We have labeled entries according to the primary task and the imaging modality.
- The SPECTER model emerges as the best model with respect to accuracy and balanced accuracy in task prediction diverse across diverse extraction strategies.
- Data and code of the paper are available at <https://github.com/zurlog/abs-embeddings-eval>.



# Future Work

## 4 Conclusions and Future Work

- To improve the annotation process of our original dataset
  - using a combination of automated tools and manual assessment.
- To collect more labeled entries in order to improve the training set sample size.
- To keep it updated.



# Abstracts Embeddings Evaluation

A Case Study of Artificial Intelligence and Medical Imaging for the COVID-19 Infection

*Thank you for listening!*

*Any questions?*